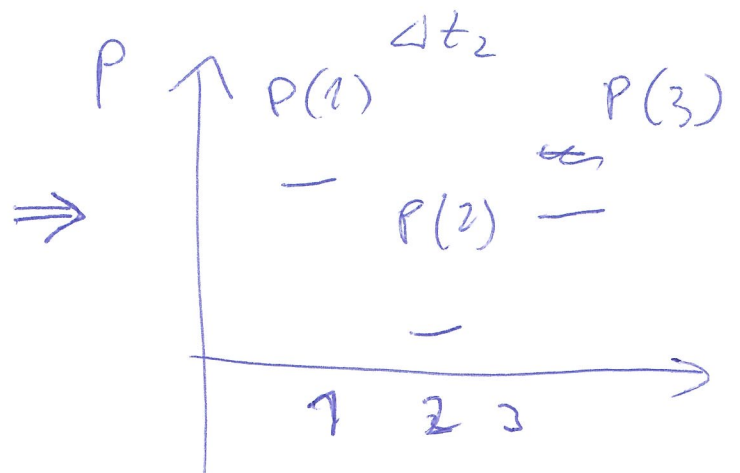
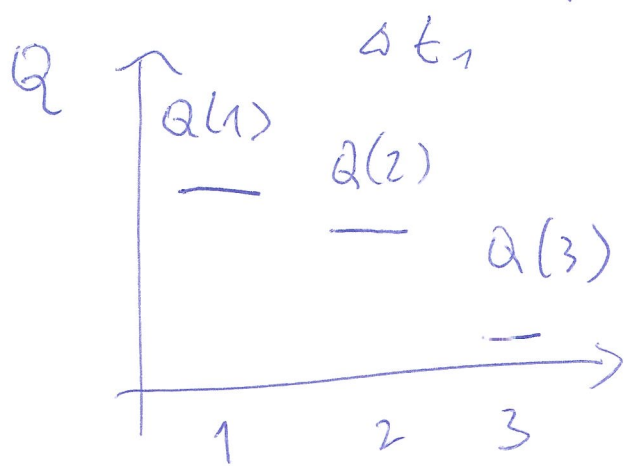


-1- KL - Divergence (Kullback-Leibler)



$P(x) \rightarrow$ new p.d.f.
 $proportion = \frac{P(x)}{Q(x)} \rightarrow$ reference p.d.f.

$KL(P \parallel Q) \rightarrow$ how P 'has changed' relative to Q
 has changed \equiv diff
 \equiv different

$$\left. \begin{array}{l} 1: \frac{P(1)}{Q(1)} = 1 \\ 2: \frac{P(2)}{Q(2)} = \frac{1}{4} \\ 3: \frac{P(3)}{Q(3)} = 4 \end{array} \right\}$$

How to use it?

\rightarrow regression function

Average (metric): $\frac{1}{2} \sum_1 \frac{P(x)}{Q(x)} =$

$$= \frac{1 + \frac{1}{4} + 4}{3} = 1.75$$

scalar
metric

-2-

-) Sometimes, asymmetry is a very defined property
-) Simple average is not
-) Mean value is not ok also because it is easily biased with outliers \rightarrow salaries are not analysed with mean (median)
-) How to balance this situation
e.g. $\frac{1}{10} \rightarrow$ balancing 10
 $p_1 \qquad p_2$

$$\frac{\sum_x p_x}{n} \rightarrow \boxed{p_1 = -p_2}$$

-) if proportion is 1 $\rightarrow p_3 = 1$, we do not want it!

$$\frac{P(x)}{Q(x)} \rightarrow \log \frac{P(x)}{Q(x)}$$

$$\frac{1}{L} \sum_l \log \frac{P(l)}{Q(l)} \rightarrow \text{'ok-ish' but still symmetrical}$$

-3-

Weighting procedure \rightarrow this is going to stress the change even more \equiv new values are more important in a sense (the opposite could also be true)

$$\sum_{\ell} P(\ell) \log \frac{P(\ell)}{Q(\ell)}$$

$w_{\ell} = P(\ell)$ \rightarrow difference with respect to Q

$$KL \Rightarrow KL(P \parallel Q) =$$

$$= \sum_{\ell} P(\ell) \log \frac{P(\ell)}{Q(\ell)}$$

$$KL(P \parallel Q) = \int_{-r}^{+r} P(x) \log \frac{P(x)}{Q(x)} dx$$

[KL is always positive or 0 if $P(x) \equiv Q(x)$]

Jensen's inequality

$$E[f(x)] \geq f(E[x])$$

$$\left(\sum_{\ell \in L} P(\ell) \log \frac{P(\ell)}{Q(\ell)} \right) \geq \log \left(\sum_{\ell \in L} P(\ell) \frac{P(\ell)}{Q(\ell)} \right)$$

-4-

$$\begin{aligned} & \log \left(\sum_{x \in L} \frac{P^2(x)}{Q(x)} \right) = \\ & = -\log \left(\sum_{x \in L} \frac{P^2(x)}{Q(x)} \right)^{-1} \geq \\ & \geq -\log \left(\frac{1}{|L|} \sum_{x \in L} \frac{P^2(x)}{Q(x)} \right) = \\ & = \log(|L|) - \log \left(\sum_{x \in L} \frac{P^2(x)}{Q(x)} \right) \end{aligned}$$

$$\text{Ex. KL}(P \parallel Q) \leq \log(|L|)$$

- 5 - Jensen-Shannon - div.

Symmetric and smoothed version of KL divergence

$P(x)$ and $Q(x)$

$$\left\{ \begin{aligned} JS(P||Q) &= \frac{1}{2} \{ KL(P||\mu) + KL(Q||\mu) \} \\ \mu &= \frac{1}{2}(P+Q) \end{aligned} \right.$$

→ average of both distributions

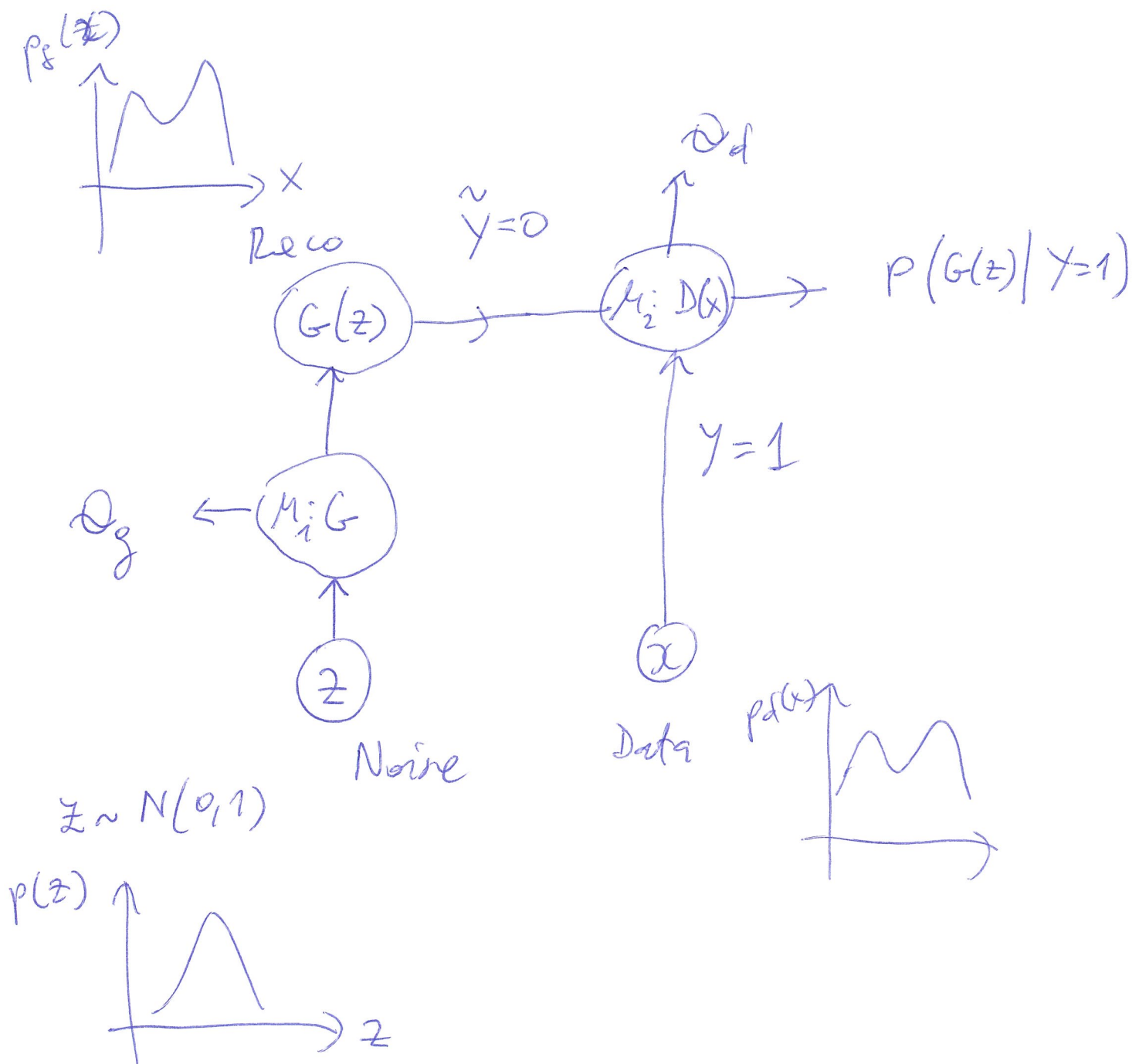
Similarity measure: small P and Q are similar, large they are different

$$\max JS(P||Q) = \log 2$$

→ useful for clustering and

→ GAN quality improvements

-6- Why GAN works at all?



$G(z) \rightarrow p_d(x)$ Universal approx. theorem (1989/1991)

G/D two mini-max players game

-7- Target function

$$\min_{G(z)} \max_{D(x)} V(G, D) = E_{x \sim p_d} [\log(D(x))] + E_{z \sim p_z} [\log(1 - D(G(z)))]$$

Formula similar to binary x-entropy

$$\rightarrow L = -y \log \tilde{y} + (1-y) \log(1-\tilde{y})$$

True data $y = 1$, $\tilde{y} = D(x)$

$$\rightarrow L_{y=1} = \log[D(x)]$$

False data $y = 0$, $\tilde{y} = D(G(z))$

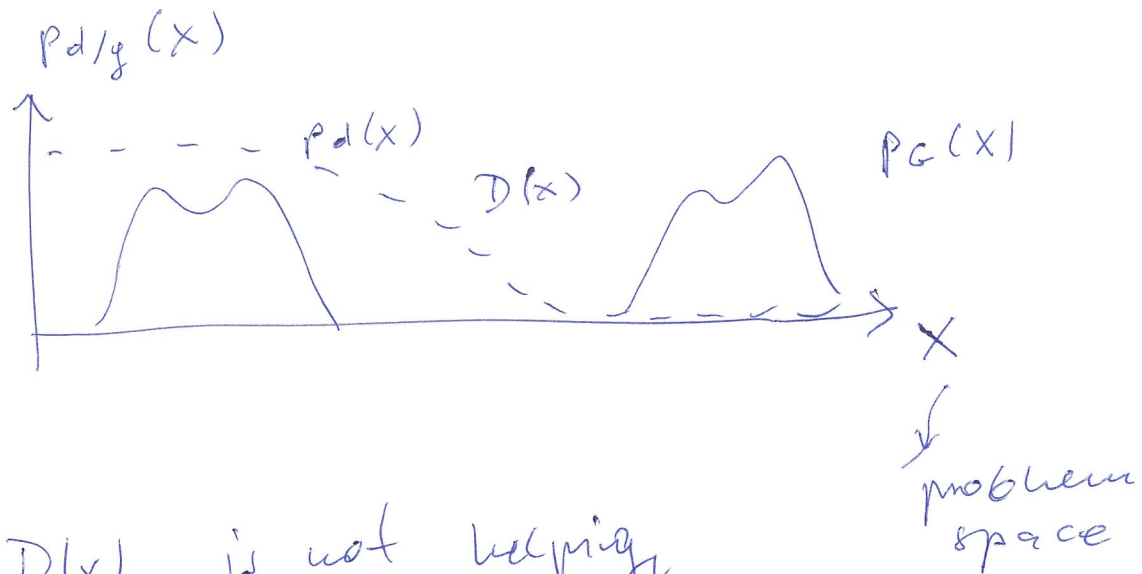
$$\rightarrow L_{y=0} = \log[1 - D(G(z))]$$

$$E[\log(D(x))] = \sum_x p_d(x_c) \log(D(x))$$

$$\int p_d(x) \log(D(x)) dx$$

$$E[\log(1 - D(G(z)))] = \sum_z p_z(z) \log(1 - D(G(z)))$$

-8- GAN's training is typically hard



$D(x)$ is not helping
in training.

→ modify D gradients

→ modify p_d and p_G distributions

→ WGAN

Typical divergences do not take into
account distance between distributions

$$p_d - p_G$$

$$W(p_d, p_G) = \inf_y E_{(x,y) \sim p(x,y)} [\|x - y\|]$$

$y \rightarrow$ that minimize the
distance

$$y_x(x) = p_d(x)$$

$$y_y(x) = p_G(x)$$

- 8 - Main loop

→ fix updates for G model

m times → inner loop updating D

•) get l instances of $y=1$ and
 l instances of $y=0$ data

•) update \mathcal{L}_D params using
grad ascent algo

$$\partial_{\mathcal{L}_D} \frac{1}{l} \{ \log(D(x)) + \log(1 - D(G(z))) \}$$

→ fix updates for D model

•) take l $y=0$ samples and
update \mathcal{L}_G using grad.
descent

$$\partial_{\mathcal{L}_G} \frac{1}{l} \{ \log(1 - D(G(z))) \}$$

Has this a chance to converge and what
the convergence mean?

Information entropy is max if all⁻¹⁻ answers are equally likely

If we go away from the equal prob. the entropy goes down (we introduced predictability)

Entropy goes down we need to ask fewer questions to guess the outcome.

Amount of information in an event \sim entropy

Information theory \rightarrow data compression, source coding

Surprise \rightarrow probability of event

$\left\{ \begin{array}{l} \text{high prob.} \rightarrow \text{low information} \\ \text{low prob.} \rightarrow \text{high information} \end{array} \right.$

$\left\{ \begin{array}{l} \text{h.p.} \rightarrow \text{h.i.} \rightarrow \text{h.s.} \\ \text{l.p.} \rightarrow \text{h.i.} \rightarrow \text{h.s.} \end{array} \right.$

$$\boxed{I(x)} = -\log_2(p(x)), \quad x \rightarrow \text{an event}$$

\downarrow unit of information is $\boxed{\text{bits}}$

$$\text{if } p(x) = 1, \quad I(x) = 0$$

Information for R.V. X with prob. distr.

$$p(X) : H(X)$$

$H(X) \rightarrow$ information entropy



The average number of bits required to represent an event drawn from the prob. distribution

$$H(X) = - \sum p \cdot \log_2(p) =$$

$$= \sum p \cdot \log_2\left(\frac{1}{p}\right)$$

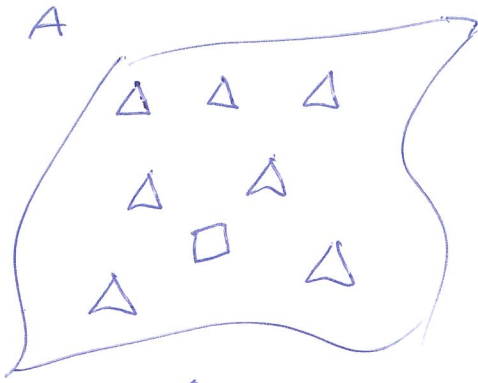
Largest entropy \rightarrow all events are equally likely

Information \equiv drop in entropy

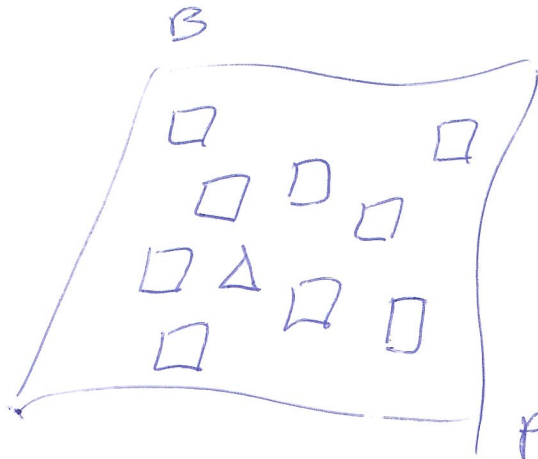
$$H(X) = - \sum_x p(x_i) \cdot \log_2(p(x_i))$$

$$H(p, q) = - \sum_x p(x_i) \cdot \log(q(x_i))$$

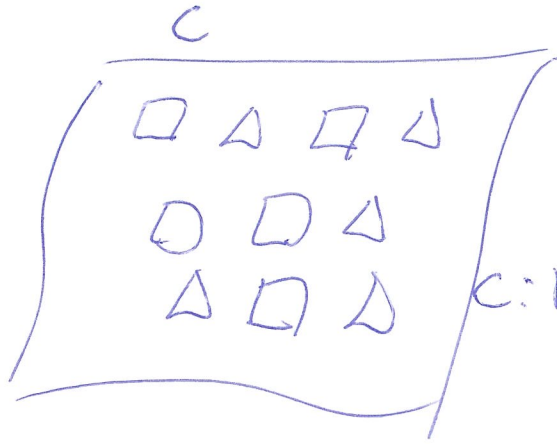
$$H(p, q) \equiv \min \text{ if } p(x_i) = q(x_i)$$



A: $P(\Delta) \uparrow$
 $P(\square) \downarrow$



B: $P(\square) \uparrow$
 $P(\Delta) \downarrow$

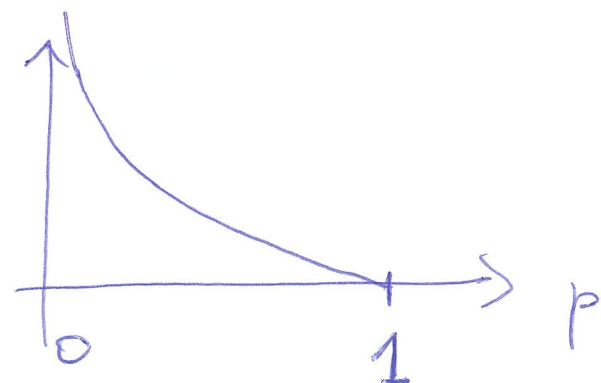


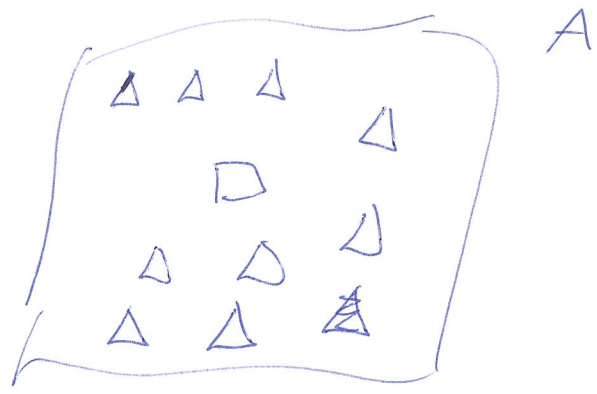
C: $P(\Delta) = P(\square)$

~~h.p.~~ A: Δ h.p., l.s.
 \square l.p., h.s.

$$S \sim \frac{1}{p}$$

$$S = \log\left(\frac{1}{p}\right)$$





$$\begin{cases} P(\Delta) = \frac{\# \Delta}{\# A_k} = \frac{9}{10} \\ P(\square) = \frac{\# \square}{\# A_k} = \frac{1}{10} \end{cases}$$

Say there is no \square , $P(\square) = 0$

$$P(\Delta) = 1, \quad S = \frac{1}{p} = \frac{1}{1} = \underline{\underline{1}}$$

Not good, we would like to get surprise = 0

$$S = \log_2 \left(\frac{1}{p} \right) = -\log_2(p)$$

$P(\Delta) \rightarrow$ getting Δ is no surprise

$P(\square) \rightarrow$ never happens, never define in this case

$$P(\Delta) = 0.8, \quad P(\square) = 0.1$$

$$S_{\Delta} = \log_2(1) - \log_2(0.8) = 0.15$$

$$S_{\square} = \log_2(1) - \log_2(0.1) = 3.32$$

$$S_{\square} > S_{\Delta}$$

Any sequence :

$$S_1 = \Delta, \Delta, \square, \Delta, \quad P(\Delta, \Delta, \square, \Delta) = 0.8 \times 0.8 \times 0.1 \times 0.8$$

$$\begin{aligned} S_1 &= -\log_2(P(S_1)) = -[3 \cdot \log_2(0.8) + \log_2(0.1)] = \\ &= 3 \cdot 0.15 + 3.32 = 3.81 \text{ [bits]} \end{aligned}$$

This works for any sequence

For instance 100 draws with return

$$[0.8 \times 100] \times 0.15 + [0.1 \times 100] \times 3.32 = 46.7$$

prob. of
 Δ in 100
draws \approx

surprise
of obs. Δ

expected
of Δ in
100 draws

\approx expected
number of Δ

$$\frac{1}{100} \downarrow$$

\equiv entropy

expected value of surprise/information is entropy

$$H(A) = 0.47$$

For the fair coin

$$H(f.c.) = \left\{ \underbrace{[0.5 \times 100]}_{\# \text{ heads}} \times 1 + \underbrace{[0.5 \times 100]}_{\# \text{ tails}} \times 1 \right\} \frac{1}{100}$$

$$= 1 \quad \text{entropy of a fair coin}$$

$$H(X) = \sum_x x \cdot p(X=x)$$

$\underbrace{\sum}_x$ $\underbrace{p(X=x)}_{\text{prob. of observing specific surprise}}$
surprise
on information

$$\begin{aligned} H(X) &= \sum_x p(x_i) \cdot \log \left(\frac{1}{p(x_i)} \right) = \\ &= - \sum_x p(x_i) \cdot \log (p(x_i)) \end{aligned}$$

$$S = k_b \cdot \ln \Omega$$

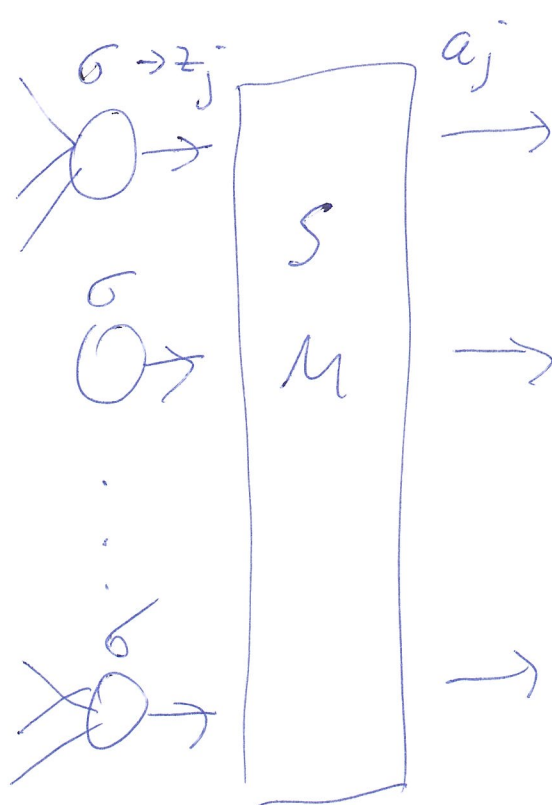
↓
disorder

$$H(x) = - \sum_x p(x_i) \log_2(p(x_i))$$

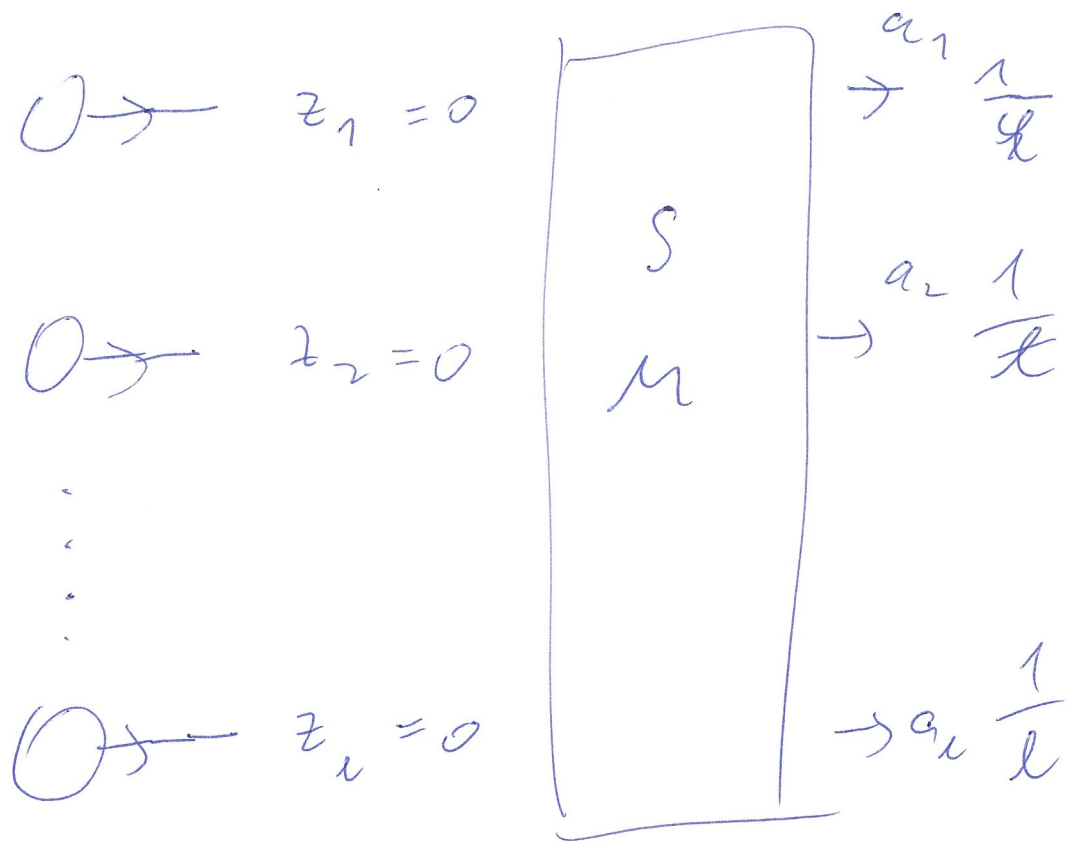
↓
uncertainty of prob. distr.

Softmax
normalisation over multiple-classes

$$\sigma_{sm}(z_j) = \frac{e^{z_j}}{\sum_l e^{z_l}}$$



a_j - probabilities



$$L_{x-e}(\tilde{y}, y) = - \sum_i y_i \log_2(\tilde{y}_i)$$

$$\tilde{y} = \begin{bmatrix} 0.1 \\ 0.173 \\ 0.56 \\ 0.167 \end{bmatrix}, \quad y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{aligned} L_{x-e} &= -0 \times \log_2(0.1) - 1 \cdot \log_2(0.173) + \\ &\quad - 0 \times \log_2(0.56) - 0 \cdot \log_2(0.167) \\ &= - \log_2(0.173) \approx 2.53 \end{aligned}$$

$$I_{MSE} = 0.68$$

$$\boxed{3\Box \quad 1\Delta}$$

$$\boxed{2\Box \quad 1\Delta}$$

Geni impurity

$$G_1 = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 \quad \left| \quad G_2 = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 \right.$$

$$= 0.38$$

$$= 0.45$$

$$\boxed{3\Box \cdot 3\Delta}$$

$$\boxed{6\Box \quad 0\Delta}$$

$$G_3 = 1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2 = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 =$$

$$= 1 - \frac{1}{2} = 0.5$$

$$G_4 = 1 - 1 = 0$$

$$G = 1 - \sum_i p_i^2$$

